

Post-course assessment

The questions are grouped into easy and moderate/hard questions. Most questions consist of both a “theory” part and an “applied” part. Choose the questions you want to work on based on your strengths and interests.

All data used in the questions on this assignment sheet can be loaded into R using the command
`load(url("http://www.stats.gla.ac.uk/~levers/aptsPCA.RData"))`

This file also contains the function `bbase` used to construct the B-spline basis in the first practical session. The model answer to the questions (except for questions 2, 4, 7 and 8 are available at the end of the document.

Easy questions

Question 1 (Kernel-density estimation). Consider the distribution represented by a density estimate \hat{f} , constructed from a sample of data $\{y_1, \dots, y_n\}$. What is the mean and variance of this distribution? What do these expressions indicate about the nature of smoothing?

Note this question does not refer to $\mathbb{E}(\hat{f}(y))$ and $\text{Var}(\hat{f}(y))$ at specific values of y , as discussed in the lectures. It refers to the mean and variance of a random variable whose density function is \hat{f} .

Question 2 (Local regression). We have defined the local mean estimator as

$$\hat{m}(x) = \frac{\sum_{k=1}^n w(x_k - x; h)y_k}{\sum_{k=1}^n w(x_k - x; h)},$$

which implies that

$$\hat{y}_i = \hat{m}(x_i) = \frac{\sum_{k=1}^n w(x_k - x_i; h)y_k}{\sum_{k=1}^n w(x_k - x_i; h)}.$$

Thus the fitted values are a linear function of the observed response and we can write $\hat{\mathbf{y}} = \mathbf{S}\mathbf{y}$.

- (a) Write down the entries of the matrix \mathbf{S} .
- (b) Consider the radio-carbon dating data (available in the data frame `radiocarbon`), in which we aim to predict the age from radio-carbon dating (`rc.age`) from the calendar age (`cal.age`).
 - (i) Construct the matrix \mathbf{S} for a Gaussian kernel with bandwidth $h = 0.05$ and use it to compute the fitted values.
 - (ii) Plot the data and add the fitted values.
 - (iii) We have seen that we can define the effective degrees of freedom as $\text{tr}(\mathbf{S})$. Compute the trace.
 - (iv) Change the value of h . How does this change the fitted function and the effective degrees of freedom?

Question 3 (Penalised least-squares for P-splines). Consider P-splines, which minimise the penalised least-squares criterion

$$\sum_{i=1}^n (y_i - m(x_i))^2 + \lambda \|\mathbf{D}\boldsymbol{\beta}\|^2 = \|\mathbf{y} - \mathbf{B}\boldsymbol{\beta}\|^2 + \lambda \|\mathbf{D}\boldsymbol{\beta}\|^2 = (\mathbf{y} - \mathbf{B}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{B}\boldsymbol{\beta}) + \lambda \boldsymbol{\beta}^\top \mathbf{D}^\top \mathbf{D} \boldsymbol{\beta},$$

where \mathbf{B} is the matrix of B-spline basis functions, \mathbf{D} is the difference matrix used in the penalty, and $\lambda \geq 0$ is the smoothing parameter.

- (a) Show, by taking the derivative with respect to β , that the minimiser of the penalised least-squares criterion is given by

$$\beta = (\mathbf{B}^\top \mathbf{B} + \lambda \mathbf{D}^\top \mathbf{D})^{-1} \mathbf{B}^\top \mathbf{y}.$$

- (b) Explain why we can rewrite the objective function as

$$\left\| \begin{pmatrix} \mathbf{B} \\ \sqrt{\lambda} \mathbf{D} \end{pmatrix} \beta - \begin{pmatrix} \mathbf{y} \\ \mathbf{0} \end{pmatrix} \right\|^2.$$

- (c) How can we exploit this to estimate β using a QR decomposition, which is numerically more stable than inverting the matrix $\mathbf{B}^\top \mathbf{B} + \lambda \mathbf{D}^\top \mathbf{D}$?

Question 4 (Clyde dissolved oxygen data). The Clyde DO data (available in the data frame `clyde`) were used in one of the practical sessions.

Subset these data to focus only on Station 10. Fit an additive model to describe the relationship between DO and the three potential explanatory variables, Doy (day of the year), Year and Salinity. Describe the results of fitting three additive terms. Now fit a model which includes interaction (a bivariate term) between Doy and Year. Is there any evidence that the seasonal effect has changed over the years at this station?

Moderate/hard questions

Question 5 (B-spline basis functions for equally-spaced knots). The B-spline basis functions are defined recursively. Given a set of l knots the B-spline basis of degree 0 is given by the functions $(B_1^0(x), \dots, B_{l-1}^0(x))$ with

$$B_j^0(x) = \begin{cases} 1 & \text{for } \kappa_j \leq x < \kappa_{j+1} \\ 0 & \text{otherwise.} \end{cases}$$

The B-spline basis of degree $r > 0$ is given by the functions $(B_1^r(x), \dots, B_{l+r-1}^r(x))$ with

$$B_j^r(x) = \frac{x - \kappa_{j-r}}{\kappa_j - \kappa_{j-r}} B_{j-1}^{r-1}(x) + \frac{\kappa_{j+1} - x}{\kappa_{j+1} - \kappa_{j+1-r}} B_j^{r-1}(x).$$

We will now turn to the important special case that the knots are equally spaced, i.e. $\kappa_{j+1} - \kappa_j = \delta$ for all j . In this case the basis functions can be computed as r -th order differences of truncated polynomials.

We start by defining the coefficients of the difference of r -th order

$$\Delta_j^r = (-1)^j \binom{r}{j} \quad \text{for } j = 0, \dots, r$$

For a first-order difference we obtain $\Delta_0^1 = 1$ and $\Delta_1^1 = -1$.

For a second order difference we obtain $\Delta_0^2 = 1$, $\Delta_1^2 = -2$ and $\Delta_2^2 = 1$.

For a third order difference we obtain $\Delta_0^3 = 1$, $\Delta_1^3 = -3$, $\Delta_2^3 = 3$ and $\Delta_3^3 = -1$.

These are also the numbers appearing inside the differencing matrix used in P-splines.

- (a) Show that $B_j^r(x) = \sum_{i=0}^{r+1} \frac{\Delta_i^{r+1}}{r! \delta^r} (x - \kappa_{j-r+i})^r$

The proof is by induction using the recursive definition of B-splines.

The following two properties of Δ_j^r will help in the proof.

- $\Delta_0^{r+1} = 1$, $\Delta_{r+1}^{r+1} = (-1)^{r+1}$ and $\Delta_j^{r+1} = \Delta_j^r - \Delta_{j-1}^r$ for $j = 1, \dots, r$.
For first-order and second-order differences this corresponds to:

$$\begin{array}{cc} 1 & -1 \\ & -1 & 1 \\ \hline 1 & -2 & 1 \end{array}$$

For second-order and third-order differences this corresponds to:

$$\begin{array}{cccc} 1 & -2 & 1 & \\ & -1 & 2 & -1 \\ \hline 1 & -3 & 3 & -1 \end{array}$$

If we define $\Delta_{-1}^r = \Delta_{r+1}^r = 0$, the formula $\Delta_j^{r+1} = \Delta_j^r - \Delta_{j-1}^r$ holds for all $j = 0, \dots, r+1$.

- $(r-i+1)\Delta_{i-1}^r = (r-i+1)(-1)^{i-1} \frac{r!}{(i-1)!(r-i+1)!} = (-1)^{i+1} \frac{r!}{(i-1)!} (r-i)! = -(-1)^i i \frac{r!}{i!(r-i)!} = -i\Delta_i^r$

- (b) Write an R script (or function) which uses the above method to generate a B-spline basis of degrees 0, 1, and 2. Can you generalise your script (or function) so that it generate a B-spline basis of any order r ?

Question 6 (Derivative estimation). In this question we will focus on estimating the derivative $m'(x)$ of the regression function. We will initially be focusing on B-splines.

- (a) Show that

$$\frac{\partial}{\partial x} B_j^r(x) = \frac{r}{\kappa_j - \kappa_{j-r}} B_{j-1}^{r-1}(x) - \frac{r}{\kappa_{j+1} - \kappa_{j+1-r}} B_j^{r-1}(x).$$

In the special case of equally-spaced knots ($\kappa_j - \kappa_{j-1} = \delta$) this becomes

$$\frac{\partial}{\partial x} B_j^r(x) = \frac{1}{\delta} B_{j-1}^{r-1}(x) - \frac{1}{\delta} B_j^{r-1}(x).$$

Hint: The proof in the general case is done by induction using the recursive definition of B-splines. In the special case of equally spaced knot, one can simply take the derivative of the formula from question 5.

- (b) Show that, in the case of equally-spaced knots,

$$\frac{\partial}{\partial x} m(x) = \sum_{j=1}^{l+r-2} B_j^{r-1}(x) \frac{\beta_{j+1} - \beta_j}{\delta},$$

i.e. we can estimate the derivative by multiplying the matrix of basis functions of degree $r-1$ with the vector $\hat{\gamma} = (\hat{\gamma}_1, \dots, \hat{\gamma}_{r+l-2})$, where $\hat{\gamma}_j = \frac{\hat{\beta}_{j+1} - \hat{\beta}_j}{\delta}$.

- (c) The data frame `follicle` contains data on the number of ovarian follicles counted from sectioned ovaries of women of various ages. It has two columns.

age	age of the women
log.count	logarithm of the follicle count

- Fit a B-spline model with a suitable number of knots to the `follicle` data.
Hint: You can use the function `bbase` and the code from the first practical session.
- Suppose we are interested in the rate by which the number of follicles is reducing. We can estimate this rate by computing the derivative of the fitted regression function. Use the formula from part (b) to compute the derivative and plot it.
- Can you construct a confidence interval for the estimated derivative? From what age onward is there a significant decrease in the number of follicles

- (d) Suppose you wanted to use a truncated power basis instead of a B-spline basis. How would you estimate the derivative?
- (e) Suppose you wanted to use a local estimate like the one studied in question 2. How would you estimate the derivative?

Question 7 (Monotonic smoothing). Consider again the radio-carbon example in which tried to relate the observed radio-carbon age to the calibrated age. It seems natural to impose the constraint that the function describing the relationship between the two is non-decreasing. In this question you will learn how this can be achieved by using equally-spaced B-splines.

- (a) Explain why the estimated regression function $\hat{m}(x) = \sum_{j=1}^{l+r-1} B_j(x)\beta_j$ is non-decreasing if $\beta_j \leq \beta_{j+1}$ for all j .

Hint: Use the derivative formula from question 6(b).

- (b) This fact can be exploited to construct a monotonic regression function. “All” we need to do is to introduce the additional constraint that $\beta_j \leq \beta_{j+1}$ for all j . These additional constraints however make finding $\hat{\beta}$ much more difficult: we have to resort to quadratic programming methods¹. We will use a simpler approach, based on modifying the penalty in a P-splines approach.² When using first-order differences to construct the penalty we use

$$\mathbf{D}_1 = \begin{pmatrix} 1 & -1 & \dots & 0 \\ \vdots & \ddots & \ddots & 0 \\ 0 & \dots & 1 & -1 \end{pmatrix},$$

i.e. the penalty becomes

$$\|\mathbf{D}_1\boldsymbol{\beta}\|^2 = \sum_{j=1}^{l+r-2} (\beta_{j+1} - \beta_j)^2.$$

In order to penalise lack of monotonicity we only want to penalise differences between the β_j if $\beta_j > \beta_{j+1}$, i.e. we would like to use the penalty

$$\sum_{j: \beta_j > \beta_{j+1}} (\beta_{j+1} - \beta_j)^2$$

This corresponds to modifying the matrix, setting all rows to zero that correspond to pairs with $\beta_j \leq \beta_{j+1}$.

Now there is of course the problem that we need the differencing matrix to estimate $\hat{\beta}$, but it, in turn, depends on β . The way around this problem is to simply iterate between these two steps, which gives the following algorithm.

1. Set $\mathbf{D}^{(1)} = \begin{pmatrix} 0 & \dots & 0 \\ \vdots & \ddots & 0 \\ 0 & \dots & 0 \end{pmatrix}$.

2. For $h = 1, 2, \dots$ until convergence ...

i. Compute $\boldsymbol{\beta}^{(h)} = (\mathbf{B}^\top \mathbf{B} + \lambda \mathbf{D}^{(h)\top} \mathbf{D}^{(h)})^{-1} \mathbf{B}^\top \mathbf{y}$ (or use a QR decomposition to compute $\boldsymbol{\beta}^{(h)}$).

ii. Set $\delta_j^{(h)} = \begin{cases} 1 & \text{if } \beta_j^{(h)} > \beta_{j+1}^{(h)} \\ 0 & \text{otherwise.} \end{cases}$

¹These are for example implemented in the R package `quadprog`.

²This idea was first suggested by Bollaerts *et al.* (British Journal of Mathematical and Statistical Psychology (2006), 59, 451–469). Bollaerts *et al.* use a monotonicity penalty in conjunction with a smoothness penalty, but for simplicity we will omit the smoothness penalty and control the smoothness by choosing a small enough number of basis functions.

$$\text{iii. Set } \mathbf{D}^{(h+1)} = \begin{pmatrix} \delta_1^{(h)} & -\delta_1^{(h)} & \dots & 0 \\ \vdots & \ddots & \ddots & 0 \\ 0 & \dots & \delta_{l+r-2}^{(h)} & -\delta_{l+r-2}^{(h)} \end{pmatrix}.$$

We will now turn to the radiocarbon data.

- (i) Fit a “classical” B-spline model with 25 equally spaced knots to the data.

Hint: You can use the function `bbase` to create the matrix \mathbf{B} .

- (ii) Implement the above algorithm to estimate a non-decreasing regression function modelling the relationship between radio-carbon age and calibrated age. Use the same basis function as in part (i).

Compare the results to the B-spline model fitted in part (i).

- (c) In this part we will return to the follicle data from question 6. All (primordial) follicles are developed before birth, so, after birth, the number of follicles must decrease over time.

When using a moderate number of knots (say 10) for modelling the log-follicle counts the estimated regression function is oscillating.

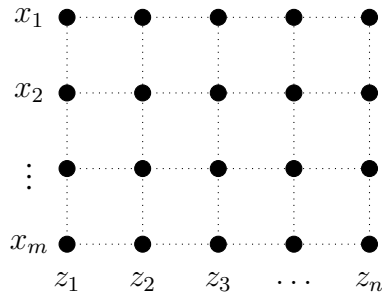
How can you use the approach set out above to estimate a *non-increasing* regression function?

Question 8 (SO₂ concentrations). The data frame `S02station` contains measurements of SO₂ from the air at a sampling station in mainland Europe. The variables are:

<code>logS02</code>	SO ₂ measurement on a log scale
<code>Year</code>	year, with a fractional component to reflect week within year
<code>Week</code>	week within the year
<code>Rain</code>	rainfall
<code>Temp</code>	temperature
<code>Humidity</code>	air humidity
<code>Flow</code>	a measurements of air flow

1. Explore the relationships between the variables, and specifically between SO₂ and potential explanatory variables, by any graphical means you consider suitable.
2. Fit an additive model and refine this into a model which you believe gives a good description of the data. (At this stage, do not worry about temporal correlation in the data.)
3. Consider a model which uses only `Year` and `Week` as explanatory variables. This model is of interest because meteorological information is not always easy to obtain, so an understanding of whether it is needed at this station may help in decisions on whether to collect it at others. Compare this reduced model with one which makes use of the meteorological information. In the comparison, interest lies particularly in any effect on the estimate of trend in SO₂ over the years.
4. Consider the additive model which contains only `Year` and `Week` as explanatory variables. Examine the residuals from this model for evidence of serial correlation. How would you adjust your model to account for this? Even if you don't do that, can you say what the effects of a suitable adjustment would be?

Question 9 (Marginal smoothing using local smoothers). Consider a bivariate smoothing problem in which the data has been collected on a regular grid, i.e. the response was observed at each combination (x_i, z_j) ($i = 1, \dots, m, j = 1, \dots, n$). The figure below illustrates this setup.



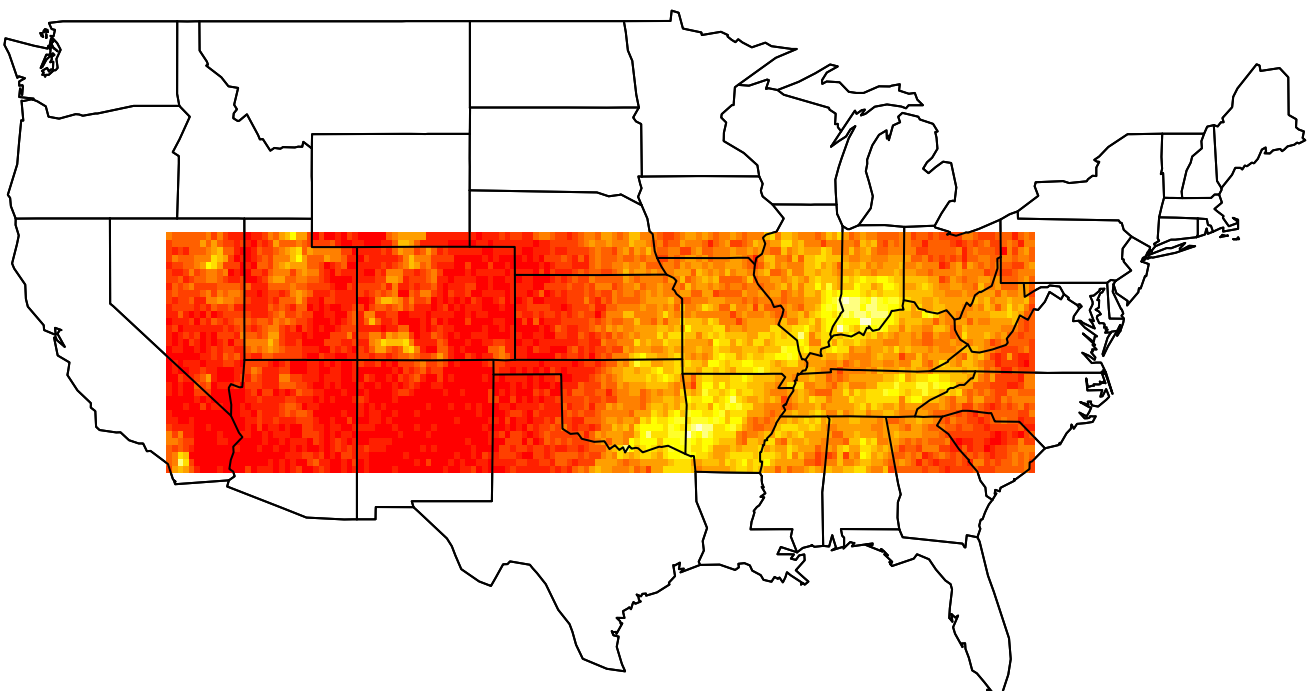
In this case it is easiest to write the observed response as a matrix $\mathbf{Y} = (y_{ij})$ with y_{ij} being the response associated with (x_i, z_j) . Alternatively we can write the observed response as a long vector, stacking the columns of \mathbf{Y} on top of each other, i.e.

$$\mathbf{y} = (y_{11}, \dots, y_{m1}, y_{12}, \dots, y_{m2}, \dots, y_{mn})^\top.$$

Suppose we wish to use a bivariate local smoother, i.e.

$$\hat{y}_{ij} = \hat{m}(x_i, z_j) = \frac{\sum_{k=1}^m \sum_{l=1}^n w(x_k - x_i; h) w(z_l - z_j; h) y_{kl}}{\sum_{k=1}^m \sum_{l=1}^n w(x_k - x_i; h) w(z_l - z_j; h)}.$$

- (a) Just like in question 2 we can write $\hat{\mathbf{y}} = \mathbf{S}\mathbf{y}$. Write down one entry of \mathbf{S} .
- (b) Show that $\mathbf{S} = \mathbf{S}^{(2)} \otimes \mathbf{S}^{(1)}$ where $\mathbf{S}^{(2)}$ is the smoothing matrix associated with a univariate smoothing problem with observed covariate values z_1, \dots, z_n and $\mathbf{S}^{(1)}$ is the smoothing matrix associated with a univariate smoothing problem with observed covariate values x_1, \dots, x_m .
Hint: Important properties of the Kronecker product (“ \otimes ”) are summarised on page 9.
- (c) Using the properties of the Kronecker product show that we can also write $\hat{\mathbf{Y}} = \mathbf{S}^{(1)}\mathbf{Y}\mathbf{S}^{(2)\top}$. What is the advantage of this representation?
Hint: Think about the dimensions of the matrices involved in the calculation.
- (d) The data frame `us.rain` contains noisy observations of the total rainfall in March/April 2006 in a rectangular area covering most of the central US (see figure below). The vectors `us.northing` and `us.easting` contain the corresponding latitudes and longitudes. Construct the smoothing matrices $\mathbf{S}^{(1)}$ and $\mathbf{S}^{(2)}$ and use the formula derived in part (c) to construct the fitted values. Use an R function like `image` to plot the resulting smoothed surface.



Question 10 (Marginal smoothing using B-splines and P-splines). Consider again the bivariate smoothing problem with gridded data set out in question 9. In this question we will consider the tensor-product-based spline approach.

(a) Explain that we can write the design matrix

$$\mathbf{B} = \begin{pmatrix} B_{11}(x_1, z_1) & \dots & B_{l_1+r-1,1}(x_1, z_1) & B_{12}(x_1, z_1) & \dots & B_{l_1+r-1,2}(x_1, z_1) & \dots & B_{l_1+r-1,l+2+r-1}(x_1, z_1) \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ B_{11}(x_m, z_1) & \dots & B_{l_1+r-1,1}(x_m, z_1) & B_{12}(x_m, z_1) & \dots & B_{l_1+r-1,2}(x_m, z_1) & \dots & B_{l_1+r-1,l+2+r-1}(x_m, z_1) \\ B_{11}(x_1, z_2) & \dots & B_{l_1+r-1,1}(x_1, z_2) & B_{12}(x_1, z_2) & \dots & B_{l_1+r-1,2}(x_1, z_2) & \dots & B_{l_1+r-1,l+2+r-1}(x_1, z_2) \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ B_{11}(x_m, z_n) & \dots & B_{l_1+r-1,1}(x_m, z_n) & B_{12}(x_m, z_n) & \dots & B_{l_1+r-1,2}(x_m, z_n) & \dots & B_{l_1+r-1,l+2+r-1}(x_m, z_n) \end{pmatrix}$$

of the bivariate tensor-product-splines as $\mathbf{B} = \mathbf{B}^{(2)} \otimes \mathbf{B}^{(1)}$, where $\mathbf{B}^{(1)}$ is a univariate B-spline basis on the x_i , $\mathbf{B}^{(2)}$ is a univariate B-spline basis on the z_j , and $B_{i,j}(x, z) = B_i^{(1)}(x)B_j^{(2)}(z)$.

(b) Suppose \mathbf{S} is the bivariate smoothing matrix, i.e.

$$\mathbf{S} = \mathbf{B}(\mathbf{B}^\top \mathbf{B})^{-1} \mathbf{B}^\top$$

and $\mathbf{S}^{(1)}$ and $\mathbf{S}^{(2)}$ are the corresponding univariate smoothing matrices, i.e.

$$\mathbf{S}^{(1)} = \mathbf{B}^{(1)}(\mathbf{B}^{(1)\top} \mathbf{B}^{(1)})^{-1} \mathbf{B}^{(1)\top} \quad \mathbf{S}^{(2)} = \mathbf{B}^{(2)}(\mathbf{B}^{(2)\top} \mathbf{B}^{(2)})^{-1} \mathbf{B}^{(2)\top}.$$

Show that $\mathbf{S} = \mathbf{S}^{(2)} \otimes \mathbf{S}^{(1)}$ and that we can write $\hat{\mathbf{y}} = \mathbf{S}\mathbf{y}$ as $\hat{\mathbf{Y}} = \mathbf{S}^{(1)}\mathbf{Y}\mathbf{S}^{(2)\top}$, just like in question 9.

(c) Is the above formula also valid if P-splines are used?

(d) Use the formula derived in part (c) to analyse the US rainfall data from question 9(d).

Question 11 (Gaussian processes).

(a) Suppose, in a Gaussian process model, we want to compute predictions $\hat{\mathbf{y}}$ for the observations we have observed data for, i.e. $\hat{y}_i = \hat{f}(\mathbf{x}_i)$. Show that these predictions are given by

$$\hat{\mathbf{y}} = \mathbf{K}(\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{y},$$

where $\mathbf{K} = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \dots & k(\mathbf{x}_1, \mathbf{x}_n) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}_n, \mathbf{x}_1) & \dots & k(\mathbf{x}_n, \mathbf{x}_n) \end{bmatrix}$ with $k(\mathbf{x}_i, \mathbf{x}_j) = \text{Cov}(Y_i, Y_j)$.

Hint: The notes contain the formula for $\hat{y}_0 = \hat{f}(\mathbf{x}_0) = \mathbb{E}(y_0|\mathbf{y})$. In this part we have to consider the case that $\mathbf{x}_0 = \mathbf{x}_i$ for $i = 1, \dots, n$, i.e. we want to evaluate the Gaussian process at the i -th observation.

(b) Show that using the spectral decomposition $\mathbf{K} = \mathbf{\Gamma} \mathbf{\Lambda} \mathbf{\Gamma}^\top$ we can write

$$\hat{\mathbf{y}} = \mathbf{\Gamma} \text{diag} \left(\frac{\lambda_i}{\lambda_i + \sigma^2} \right) \mathbf{\Gamma}^\top \mathbf{y}.$$

(c) Suppose that the data is on a regular grid like in question 9 and suppose that we use a separable covariance/kernel function, i.e.

$$k((x_i, z_j), (x_k, z_l)) = k^{(1)}(x_i, x_k)k^{(2)}(z_j, z_l)$$

(i) Show that we can write $\mathbf{K} = \mathbf{K}^{(2)} \otimes \mathbf{K}^{(1)}$, where

$$\mathbf{K}^{(1)} = \begin{bmatrix} k^{(1)}(x_1, x_1) & \dots & k^{(1)}(x_1, x_m) \\ \vdots & \ddots & \vdots \\ k^{(1)}(x_m, x_1) & \dots & k^{(1)}(x_m, x_m) \end{bmatrix} \quad \mathbf{K}^{(2)} = \begin{bmatrix} k^{(2)}(z_1, z_1) & \dots & k^{(2)}(z_1, z_n) \\ \vdots & \ddots & \vdots \\ k^{(2)}(z_n, z_1) & \dots & k^{(2)}(z_n, z_n) \end{bmatrix}$$

- (ii) Using the properties of the Kronecker product show that in this case we can compute the entries of the vector $\hat{\mathbf{y}}$ as

$$\hat{\mathbf{Y}} = \mathbf{\Gamma}^{(1)} \left(\mathbf{D} \odot \left(\mathbf{\Gamma}^{(1)\top} \mathbf{Y} \mathbf{\Gamma}^{(2)} \right) \right) \mathbf{\Gamma}^{(2)\top},$$

with $\mathbf{K}^{(1)} = \mathbf{\Gamma}^{(1)} \mathbf{\Lambda}^{(1)} \mathbf{\Gamma}^{(1)\top}$ and $\mathbf{K}^{(2)} = \mathbf{\Gamma}^{(2)} \mathbf{\Lambda}^{(2)} \mathbf{\Gamma}^{(2)\top}$ being the spectral decompositions of $\mathbf{K}^{(1)}$ and $\mathbf{K}^{(2)}$ and $D_{ij} = \frac{\lambda_i^{(1)} \lambda_j^{(2)}}{\lambda_i^{(1)} \lambda_j^{(2)} + \sigma^2}$. \odot denotes element-wise multiplication of matrices (“*” instead of “%*%” in R-speak).

The key advantage of this formula is that it only involves working with matrices of dimensions $m \times n$, $m \times m$ and $n \times n$, whereas a naïve implementation would require working with matrices of dimension $(mn) \times (mn)$, which are a lot bigger.

- (d) (i) Compute $\hat{\mathbf{y}}$ for the US rain data by using a kernel/covariance function of your choice. What effect does changing the hyperparameters have?

Hint: Reasonable values for the hyperparameters are $\tau^2 \approx 6000$, $\sigma^2 \approx 500$ and when using a Matérn kernel (function `matern` in `geoR`) $\phi \approx 0.7$ and $\kappa = 1.5$. Alternatively you can use a Gaussian kernel; `geoR`'s $\phi = 0.7$ corresponds to $\rho = 2$ using the notation from the lectures.

- (ii) Can you implement the more efficient formula used in part (ii) of (c)? (You need to assume a separable covariance structure for this approach to apply).

- (iii) Can you obtain the same results using the `krige.conv` function from the package `geoR`?

Hint: Estimate the hyperparameters first using `likfit` or `vario.est`.

The Kronecker Product

Given a $m \times n$ matrix \mathbf{C} and a $p \times q$ matrix \mathbf{D} the *Kronecker product* of the two matrices is defined as the following $mp \times nq$ matrix:

$$\mathbf{C} \otimes \mathbf{D} = \begin{bmatrix} c_{11}\mathbf{D} & \cdots & c_{1n}\mathbf{D} \\ \vdots & \ddots & \vdots \\ c_{m1}\mathbf{D} & \cdots & c_{mn}\mathbf{D} \end{bmatrix} = \begin{bmatrix} c_{11}d_{11} & c_{11}d_{12} & \cdots & c_{11}d_{1q} & \cdots & c_{1n}d_{11} & c_{1n}d_{12} & \cdots & c_{1n}d_{1q} \\ c_{11}d_{21} & c_{11}d_{22} & \cdots & c_{11}d_{2q} & \cdots & c_{1n}d_{21} & c_{1n}d_{22} & \cdots & c_{1n}d_{2q} \\ \vdots & \vdots & \ddots & \vdots & & \vdots & \vdots & \ddots & \vdots \\ c_{11}d_{p1} & c_{11}d_{p2} & \cdots & c_{11}d_{pq} & \cdots & c_{1n}d_{p1} & c_{1n}d_{p2} & \cdots & c_{1n}d_{pq} \\ \vdots & \vdots & & \vdots & \ddots & \vdots & \vdots & & \vdots \\ c_{m1}d_{11} & c_{m1}d_{12} & \cdots & c_{m1}d_{1q} & \cdots & c_{mn}d_{11} & c_{mn}d_{12} & \cdots & c_{mn}d_{1q} \\ c_{m1}d_{21} & c_{m1}d_{22} & \cdots & c_{m1}d_{2q} & \cdots & c_{mn}d_{21} & c_{mn}d_{22} & \cdots & c_{mn}d_{2q} \\ \vdots & \vdots & \ddots & \vdots & & \vdots & \vdots & \ddots & \vdots \\ c_{m1}d_{p1} & c_{m1}d_{p2} & \cdots & c_{m1}d_{pq} & \cdots & c_{mn}d_{p1} & c_{mn}d_{p2} & \cdots & c_{mn}d_{pq} \end{bmatrix}$$

For this assignment you will need the following properties of the Kronecker product.

- The Kronecker product is, just like the standard matrix product, *not* commutative, i.e. $\mathbf{C} \otimes \mathbf{D} \neq \mathbf{D} \otimes \mathbf{C}$. It is however associative, i.e. $(\mathbf{C} \otimes \mathbf{D}) \otimes \mathbf{E} = \mathbf{C} \otimes (\mathbf{D} \otimes \mathbf{E})$.
- $\mathbf{C} \otimes (\mathbf{D} + \mathbf{E}) = \mathbf{C} \otimes \mathbf{D} + \mathbf{C} \otimes \mathbf{E}$ and $(\mathbf{C} + \mathbf{D}) \otimes \mathbf{E} = \mathbf{C} \otimes \mathbf{E} + \mathbf{D} \otimes \mathbf{E}$.
- $(\mathbf{C} \otimes \mathbf{D})(\mathbf{E} \otimes \mathbf{F}) = (\mathbf{CE}) \otimes (\mathbf{DF})$.
- $(\mathbf{C} \otimes \mathbf{D})^{-1} = \mathbf{C}^{-1} \otimes \mathbf{D}^{-1}$.
- $(\mathbf{C} \otimes \mathbf{D})^\top = \mathbf{C}^\top \otimes \mathbf{D}^\top$.
- $(\mathbf{A} \otimes \mathbf{B})\mathbf{c} = \mathbf{d}$ if and only if $\mathbf{BCA}^\top = \mathbf{D}$, where the vectors \mathbf{c} (and \mathbf{d}) simply consist of the columns of \mathbf{C} (and \mathbf{D}) stacked on top of each other, i.e.

$$\mathbf{C} = \begin{pmatrix} c_{11} & \cdots & c_{1n} \\ \vdots & \ddots & \vdots \\ c_{m1} & \cdots & c_{mn} \end{pmatrix} \quad \mathbf{c} = (c_{11}, \dots, c_{m1}, c_{12}, \dots, c_{m2}, \dots, c_{mn})^\top$$

$$\mathbf{D} = \begin{pmatrix} d_{11} & \cdots & d_{1q} \\ \vdots & \ddots & \vdots \\ d_{p1} & \cdots & d_{pq} \end{pmatrix} \quad \mathbf{d} = (d_{11}, \dots, d_{p1}, d_{12}, \dots, d_{p2}, \dots, d_{pq})^\top$$

\mathbf{A} is a $q \times n$ matrix and \mathbf{B} is a $p \times m$ matrix, thus \mathbf{C} is a $m \times n$ matrix and \mathbf{D} is a $p \times q$ matrix.

In R this can be verified as follows.

```
R1 # Create example matrices
R2 A <- matrix(rnorm(6), ncol=3)
R3 B <- matrix(rnorm(8), ncol=2)
R4 C <- matrix(rnorm(6), ncol=3)
R5 c <- as.vector(C)
R6 # Calculate d
R7 kronecker(A,B)%*%c
R8 # Calculate D
R9 B%*%C%*%t(A)
```

- Consider two square matrices \mathbf{A} and \mathbf{B} with spectral decompositions $\mathbf{A} = \mathbf{\Gamma}_A \mathbf{\Lambda}_A \mathbf{\Gamma}_A^\top$ and $\mathbf{B} = \mathbf{\Gamma}_B \mathbf{\Lambda}_B \mathbf{\Gamma}_B^\top$, i.e. the $\mathbf{\Gamma}$ matrices contain the eigenvectors and $\mathbf{\Lambda}$ is a diagonal matrix with the corresponding eigenvalues. Then $\mathbf{A} \otimes \mathbf{B}$ has the spectral decomposition

$$\mathbf{A} \otimes \mathbf{B} = (\mathbf{\Gamma}_A \otimes \mathbf{\Gamma}_B) (\mathbf{\Lambda}_A \otimes \mathbf{\Lambda}_B) (\mathbf{\Gamma}_A \otimes \mathbf{\Gamma}_B)^\top.$$

Post-course assessment – Model answers

1. Mean: $\int \hat{f}(y)ydy = \frac{1}{n} \sum_{i=1}^n \int y \frac{1}{h} w\left(\frac{y-y_i}{h}\right) dy = \frac{1}{n} \sum_{i=1}^n y_i = \bar{y}$

$$\int \hat{f}(y)y^2dy = \frac{1}{n} \sum_{i=1}^n \int y^2 \frac{1}{h} w\left(\frac{y-y_i}{h}\right) dy = \frac{1}{n} \sum_{i=1}^n (h^2 + y_i^2)$$

$$\text{Variance: } h^2 + \frac{1}{n} \sum_{i=1}^n y_i^2 - \bar{y}^2 = h^2 + \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2$$

So, the mean of the distribution is the sample mean of the data, while the variance is the sample variance (using the divisor n) inflated by h^2 . This ‘variance inflation’ quantifies the effect of smoothing.

2. *Model answers not being made available.*

3. (a) We have that

$$\begin{aligned} & \frac{\partial}{\partial \boldsymbol{\beta}} (\mathbf{y} - \mathbf{B}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{B}\boldsymbol{\beta}) + \lambda \boldsymbol{\beta}^\top \mathbf{D}^\top \mathbf{D} \boldsymbol{\beta} \\ &= \frac{\partial}{\partial \boldsymbol{\beta}} \mathbf{y}^\top \mathbf{y} - 2\boldsymbol{\beta}^\top \mathbf{B}^\top \mathbf{y} + \boldsymbol{\beta}^\top \mathbf{B}^\top \mathbf{B} \boldsymbol{\beta} + \lambda \boldsymbol{\beta}^\top \mathbf{D}^\top \mathbf{D} \boldsymbol{\beta} \\ &= -2\mathbf{B}^\top \mathbf{y} + \mathbf{B}^\top \mathbf{B} \boldsymbol{\beta} + \lambda \mathbf{D}^\top \mathbf{D} \boldsymbol{\beta} \end{aligned}$$

Setting this derivative to zero yields the equation

$$(\mathbf{B}^\top \mathbf{B} + \lambda \mathbf{D}^\top \mathbf{D}) \boldsymbol{\beta} = \mathbf{B}^\top \mathbf{y}$$

which has

$$\hat{\boldsymbol{\beta}} = (\mathbf{B}^\top \mathbf{B} + \lambda \mathbf{D}^\top \mathbf{D})^{-1} \mathbf{B}^\top \mathbf{y}$$

as solution.

(b) Using the formula that $\begin{pmatrix} \mathbf{A} \\ \mathbf{B} \end{pmatrix}^\top \begin{pmatrix} \mathbf{C} \\ \mathbf{D} \end{pmatrix} = \mathbf{A}^\top \mathbf{C} + \mathbf{B}^\top \mathbf{D}$ we have that

$$\begin{aligned} & (\mathbf{y} - \mathbf{B}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{B}\boldsymbol{\beta}) + \lambda \boldsymbol{\beta}^\top \mathbf{D}^\top \mathbf{D} \boldsymbol{\beta} \\ &= (\mathbf{y} - \mathbf{B}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{B}\boldsymbol{\beta}) + (\mathbf{0} - \sqrt{\lambda} \mathbf{D} \boldsymbol{\beta})^\top (\mathbf{0} - \sqrt{\lambda} \mathbf{D} \boldsymbol{\beta}) \\ &= \left(\begin{pmatrix} \mathbf{y} \\ \mathbf{0} \end{pmatrix} - \begin{pmatrix} \mathbf{B} \\ \sqrt{\lambda} \mathbf{D} \end{pmatrix} \boldsymbol{\beta} \right)^\top \left(\begin{pmatrix} \mathbf{y} \\ \mathbf{0} \end{pmatrix} - \begin{pmatrix} \mathbf{B} \\ \sqrt{\lambda} \mathbf{D} \end{pmatrix} \boldsymbol{\beta} \right) \\ &= \left\| \begin{pmatrix} \mathbf{y} \\ \mathbf{0} \end{pmatrix} - \begin{pmatrix} \mathbf{B} \\ \sqrt{\lambda} \mathbf{D} \end{pmatrix} \boldsymbol{\beta} \right\|^2 \end{aligned}$$

(c) The QR decomposition of $\tilde{\mathbf{B}} = \begin{pmatrix} \mathbf{B} \\ \sqrt{\lambda} \mathbf{D} \end{pmatrix}$ is $\tilde{\mathbf{B}} = (\mathbf{Q}_1 \ \mathbf{Q}_2) \begin{pmatrix} \mathbf{R} \\ \mathbf{0} \end{pmatrix} = \mathbf{Q}_1 \mathbf{R}$. $\mathbf{Q} = (\mathbf{Q}_1 \ \mathbf{Q}_2)$ is an orthogonal matrix and \mathbf{R} is upper triangular. With $\tilde{\mathbf{y}} = \begin{pmatrix} \mathbf{y} \\ \mathbf{0} \end{pmatrix}$ we have that

$$\begin{aligned} \|\tilde{\mathbf{y}} - \tilde{\mathbf{B}}\boldsymbol{\beta}\|^2 &= \|\tilde{\mathbf{y}} - \mathbf{Q}_1 \mathbf{R} \boldsymbol{\beta}\|^2 = \|\mathbf{Q}^\top (\tilde{\mathbf{y}} - \mathbf{Q}_1 \mathbf{R} \boldsymbol{\beta})\|^2 \\ &= \left\| \begin{pmatrix} \mathbf{Q}_1^\top \tilde{\mathbf{y}} \\ \mathbf{Q}_2^\top \tilde{\mathbf{y}} \end{pmatrix} - \begin{pmatrix} \mathbf{R} \boldsymbol{\beta} \\ \mathbf{0} \end{pmatrix} \right\|^2 = \|\mathbf{Q}_1^\top \tilde{\mathbf{y}} - \mathbf{R} \boldsymbol{\beta}\|^2 + \|\mathbf{Q}_2^\top \tilde{\mathbf{y}}\|^2 \end{aligned}$$

We have used that $\mathbf{Q}_1^\top \mathbf{Q}_1 = \mathbf{I}$ and $\mathbf{Q}_1^\top \mathbf{Q}_2 = \mathbf{0}$.

$\|\mathbf{Q}_2^\top \tilde{\mathbf{y}}\|^2$ does not depend on $\boldsymbol{\beta}$ and $\|\mathbf{Q}_1^\top \tilde{\mathbf{y}} - \mathbf{R} \boldsymbol{\beta}\|^2$ can (if $\tilde{\mathbf{B}}$ is of full column rank) be made exactly zero by solving $\mathbf{R} \boldsymbol{\beta} = \mathbf{Q}_1^\top \tilde{\mathbf{y}}$, which can, due to \mathbf{R} being upper-triangular, be performed very efficiently.


```

R 12 # Create knots
R 13 knots <- seq(min(x)-delta, max(x)+delta, delta)
R 14 # Create truncated polynomials
R 15 B <- outer(x, knots, function(x,y) pmax(x-y,0))
R 16 # Create B-spline matrix
R 17 B <- (B[,1:(ncol(B)-2)] - 2*B[,2:(ncol(B)-1)] + B[, 3:ncol(B)])/delta
R 18 # Plot basis functions
R 19 matplot(x, B, type="l")

```

Degree $r = 2$:

```

R 20 # Create knots
R 21 knots <- seq(min(x)-2*delta, max(x)+2*delta, delta)
R 22 # Create truncated polynomials
R 23 B <- outer(x, knots, function(x,y) pmax(x-y,0))^2
R 24 # Create B-spline matrix
R 25 B <- (B[,1:(ncol(B)-3)] - 3*B[,2:(ncol(B)-2)] + 3*B[,3:(ncol(B)-1)] - B[, 4:ncol(B)])
R 26 / (2*delta^2)
R 27 # Plot basis functions
R 28 matplot(x, B, type="l")

```

Any degree r :

```

R 29 # Set example degree
R 30 r <- 3
R 31 # Create knots
R 32 knots <- seq(min(x)-r*delta, max(x)+r*delta, delta)
R 33 # Create truncated polynomials
R 34 B <- outer(x, knots, function(x,y) pmax(x-y,0))^r
R 35 # Create B-spline matrix
R 36 B <- B%*%t(diff(diag(ncol(B)), diff=r+1))/(delta^r*factorial(r))
R 37 # Plot basis functions
R 38 matplot(x, B, type="l")

```

6. (a) The proof is by induction.

$r = 1$: For $r = 1$ we have that

$$B_j^r(x) = \begin{cases} \frac{x-\kappa_{j-1}}{\kappa_j-\kappa_{j-1}} & \text{for } \kappa_{j-1} \leq x < \kappa_j \\ \frac{\kappa_j-x}{\kappa_{j+1}-\kappa_j} & \text{for } \kappa_j \leq x < \kappa_{j+1} \\ 0 & \text{otherwise,} \end{cases}$$

thus

$$\begin{aligned} \frac{\partial}{\partial x} B_j^r(x) &= \begin{cases} \frac{1}{\kappa_j-\kappa_{j-1}} & \text{for } \kappa_{j-1} \leq x < \kappa_j \\ \frac{-1}{\kappa_{j+1}-\kappa_j} & \text{for } \kappa_j \leq x < \kappa_{j+1} \\ 0 & \text{otherwise,} \end{cases} \\ &= \frac{1}{\kappa_j - \kappa_{j-1}} B_{j-1}^{r-1}(x) - \frac{1}{\kappa_{j+1} - \kappa_j} B_j^{r-1}(x) \end{aligned}$$

$r - 1 \rightsquigarrow r$: Suppose we know that

$$\frac{\partial}{\partial x} B_j^{r-1}(x) = -\frac{r-1}{\kappa_j - \kappa_{j-r+1}} B_{j-1}^{r-2}(x) + \frac{r-1}{\kappa_{j+1} - \kappa_{j-r+2}} B_j^{r-2}(x).$$

We will take the derivative of $B_j^r(x)$.

$$\begin{aligned} \frac{\partial}{\partial x} B_j^r(x) &= \frac{1}{\kappa_j - \kappa_{j-r}} B_{j-1}^{r-1}(x) + \frac{x - \kappa_{j-r}}{\kappa_j - \kappa_{j-r}} \frac{\partial}{\partial x} B_{j-1}^{r-1}(x) \\ &\quad - \frac{1}{\kappa_{j+1} - \kappa_{j+1-r}} B_j^{r-1}(x) + \frac{\kappa_{j+1} - x}{\kappa_{j+1} - \kappa_{j+1-r}} \frac{\partial}{\partial x} B_j^{r-1}(x) \end{aligned} \quad (1)$$

We will now rewrite the right-most term in each of the last two lines of the above equation.

$$\begin{aligned} \frac{x - \kappa_{j-r}}{\kappa_j - \kappa_{j-r}} \frac{\partial}{\partial x} B_{j-1}^{r-1}(x) &= \frac{x - \kappa_{j-r}}{\kappa_j - \kappa_{j-r}} \left(\frac{r-1}{\kappa_{j-1} - \kappa_{j-r}} B_{j-2}^{r-2}(x) - \frac{r-1}{\kappa_j - \kappa_{j+1-r}} B_{j-1}^{r-2}(x) \right) \\ &= \frac{r-1}{\kappa_j - \kappa_{j-r}} \underbrace{\left(\frac{x - \kappa_{j-r}}{\kappa_{j-1} - \kappa_{j-r}} B_{j-2}^{r-2}(x) + \frac{\kappa_j - x}{\kappa_j - \kappa_{j+1-r}} B_{j-1}^{r-2}(x) \right)}_{=-B_{j-1}^{r-1}(x)} \\ &\quad - \underbrace{\left(\frac{\kappa_j - x}{\kappa_j - \kappa_{j+1-r}} B_{j-1}^{r-2}(x) - \frac{x - \kappa_{j-r}}{\kappa_j - \kappa_{j+1-r}} B_{j-1}^{r-2}(x) \right)}_{=-\frac{\kappa_j - \kappa_{j-r}}{\kappa_j - \kappa_{j+1-r}} B_{j-1}^{r-2}(x)} \\ &= \frac{r-1}{\kappa_j - \kappa_{j-r}} B_{j-1}^{r-1}(x) - \frac{r-1}{\kappa_j - \kappa_{j+1-r}} B_{j-1}^{r-2}(x) \\ \\ \frac{\kappa_{j+1} - x}{\kappa_{j+1} - \kappa_{j+1-r}} \frac{\partial}{\partial x} B_j^{r-1}(x) &= \frac{\kappa_{j+1} - x}{\kappa_{j+1} - \kappa_{j+1-r}} \left(\frac{r-1}{\kappa_j - \kappa_{j+1-r}} B_{j-1}^{r-2}(x) - \frac{r-1}{\kappa_{j+1} - \kappa_{j+2-r}} B_j^{r-2}(x) \right) \\ &= \frac{r-1}{\kappa_{j+1} - \kappa_{j+1-r}} \underbrace{\left(\frac{\kappa_{j+1} - x}{\kappa_j - \kappa_{j+1-r}} B_{j-1}^{r-2}(x) + \frac{x - \kappa_{j-r+1}}{\kappa_j - \kappa_{j-r+1}} B_{j-1}^{r-2}(x) \right)}_{=\frac{\kappa_{j+1} - \kappa_{j-r+1}}{\kappa_j - \kappa_{j-r+1}} B_{j-1}^{r-2}(x)} \\ &\quad - \underbrace{\left(\frac{x - \kappa_{j-r+1}}{\kappa_j - \kappa_{j-r+1}} B_{j-1}^{r-2}(x) - \frac{\kappa_{j+1} - x}{\kappa_{j+1} - \kappa_{j+2-r}} B_j^{r-2}(x) \right)}_{=-B_j^{r-1}(x)} \\ &= \frac{r-1}{\kappa_j - \kappa_{j-r+1}} B_{j-1}^{r-2}(x) - \frac{r-1}{\kappa_{j+1} - \kappa_{j+1-r}} B_j^{r-1}(x) \end{aligned}$$

In both cases the second and third term inserted in the second step cancel each other out (second term is minus the third term). Plugging these two results into equation (1) yields

$$\begin{aligned} \frac{\partial}{\partial x} B_j^r(x) &= \frac{1}{\kappa_j - \kappa_{j-r}} B_{j-1}^{r-1}(x) + \frac{r-1}{\kappa_j - \kappa_{j-r}} B_{j-1}^{r-1}(x) - \frac{r-1}{\kappa_j - \kappa_{j+1-r}} B_{j-1}^{r-2}(x) \\ &\quad - \frac{1}{\kappa_{j+1} - \kappa_{j+1-r}} B_j^{r-1}(x) + \frac{r-1}{\kappa_j - \kappa_{j-r+1}} B_{j-1}^{r-2}(x) - \frac{r-1}{\kappa_{j+1} - \kappa_{j+1-r}} B_j^{r-1}(x) \\ &= \frac{r}{\kappa_j - \kappa_{j-r}} B_{j-1}^{r-1}(x) - \frac{r}{\kappa_{j+1} - \kappa_{j+1-r}} B_j^{r-1}(x) \end{aligned}$$

The proof in the special case of equally-spaced knots is much simpler and does not need to be per-

formed using induction.

$$\begin{aligned}
\frac{\partial}{\partial x} B_j^r(x) &= \frac{\partial}{\partial x} \sum_{i=0}^{r+1} \frac{\Delta_i^{r+1}}{r! \delta^r} (x - \kappa_{j-r+i})_+^r \\
&= \sum_{i=0}^{r+1} \frac{\overbrace{\Delta_i^{r+1}}^{\Delta_i - \Delta_{i-1}}}{(r-1)! \delta^r} (x - \kappa_{j-r+i})_+^{r-1} \\
&= \frac{1}{\delta} \sum_{i=0}^r \frac{\Delta_i^r}{(r-1)! \delta^{r-1}} (x - \kappa_{j-r+i})_+^{r-1} - \frac{1}{\delta} \sum_{i=1}^{r+1} \frac{\Delta_{i-1}^r}{(r-1)! \delta^{r-1}} (x - \kappa_{j-r+i})_+^{r-1} \\
&= \frac{1}{\delta} \underbrace{\sum_{i=0}^r \frac{\Delta_i^r}{(r-1)! \delta^{r-1}} (x - \kappa_{j-1-r+1+i})_+^{r-1}}_{=B_{j-1}^{r-1}(x)} - \frac{1}{\delta} \underbrace{\sum_{i=0}^r \frac{\Delta_i^r}{(r-1)! \delta^{r-1}} (x - \kappa_{j-r+1+i})_+^{r-1}}_{=B_j^{r-1}(x)}
\end{aligned}$$

(b) Using the derivatives of the basis functions we have just derived

$$\begin{aligned}
\frac{\partial}{\partial x} m(x) &= \sum_{j=1}^{l+r-1} \frac{\partial}{\partial x} B_j^r(x) \beta_j \\
&= \sum_{j=1}^{l+r-1} (B_{j-1}^{r-1}(x) - B_j^{r-1}(x)) \frac{\beta_j}{\delta} \\
&= \sum_{j=1}^{l+r-2} B_j^{r-1}(x) \frac{\beta_{j+1} - \beta_j}{\delta}
\end{aligned}$$

In the last step we have exploited that within the range of the data $B_0^{r-1}(x) = B_r^{l+r-1}(x) = 0$.

(c) (parts (i) to (iii))

```

R1 # Set x to covariate and y to response
R2 x <- follicle$age
R3 y <- follicle$log.count
R4
R5 # Create basis functions
R6 B <- bbase(x, deg=3, n.knots=4)
R7 # Estimate coefficients
R8 beta <- qr.coef(qr(B), y)
R9
R10 # Plot fitted function
R11 plot(x, y, xlab="Age", ylab="log(Count)")
R12 xx<-seq(min(x), max(x), len=1e3)
R13 BB <- bbase(xx, deg=3, n.knots=4)
R14 lines(xx, BB%*%beta)
R15
R16 # Estimated derivative
R17 BBd <-bbase(xx, deg=2, n.knots=4)
R18 delta <- (max(x)-min(x))/(4-1)
R19 gamma <- diff(beta)/delta
R20 deriv <- BBd%*%gamma
R21 plot(xx, deriv, type="l", xlab="Age", ylab="Derivative")

```

The estimated derivative is highly sensitive to the choice of number of basis functions.

Just like in the classical linear model we have that

$$\text{Var}(\hat{\beta}) = \sigma^2 * (\mathbf{B}'\mathbf{B})^{-1}$$

The estimated derivative is $\mathbf{B}^{r-1}\hat{\gamma} = \frac{1}{\delta}\mathbf{B}^{r-1}\mathbf{D}_1\hat{\beta}$ where $\mathbf{D}_1 = \begin{pmatrix} 1 & -1 & \dots & 0 \\ \vdots & \ddots & \ddots & 0 \\ 0 & \dots & 1 & -1 \end{pmatrix}$ is the first order

differencing matrix and \mathbf{B}^{r-1} is the corresponding B-spline basis matrix of degree $r - 1$. Thus the variance of the estimated derivative is

$$\text{Var}\left(\frac{1}{\delta}\mathbf{B}^{r-1}\mathbf{D}_1\hat{\beta}\right) = \frac{\sigma^2}{\delta^2}\mathbf{B}^{r-1}\mathbf{D}_1(\mathbf{B}'\mathbf{B})^{-1}\mathbf{D}_1'\mathbf{B}^{r-1}$$

The square roots of the diagonal elements of the matrix are the standard error of the derivatives.

```
R 22 # Estimate residual sd
R 23 sigma.hat <- sqrt(1/(length(y)-ncol(B)) * sum((B%*%beta - y)^2))
R 24
R 25 # Compute standard error of derivatives
R 26 D <- diff(diag(ncol(B)))
R 27 sds <- sigma.hat/delta * sqrt(diag(BBd%*%D%*%solve(crossprod(B))%*%t(D)%*%t(BBd)))
R 28
R 29 # Plot confidence bands
R 30 lines(xx, deriv + qnorm(0.975) * sds, lty=2)
R 31 lines(xx, deriv - qnorm(0.975) * sds, lty=2)
R 32
R 33 # Add horizontal line at 0
R 34 abline(h=0, lty=3)
```

(d) The basis functions of a truncated power basis of degree r with l knots is

$$1, x, \dots, x^{r-1}, (x - \kappa_1)_+^r, \dots, (x - \kappa_{l-1})^r$$

Taking derivatives yields

$$0, 1, \dots, (r - 1)x^{r-2}, r(x - \kappa_1)_+^{r-1}, \dots, r(x - \kappa_{l-1})^{r-1}$$

Thus the derivative of the regression function can be computed as

$$m'(\mathbf{x}) = \mathbf{B}^{-1}\boldsymbol{\gamma}$$

where \mathbf{B}^{r-1} is the matrix of truncated power functions of degree $r - 1$ and $\boldsymbol{\gamma} = (\beta_1, 2\beta_2, \dots, (r - 1)\beta_{r-1}, r\beta_r, \dots, r\beta_{l+r-1})^\top$.

```
R 35 # Create basis functions (using tbase function from practical 1)
R 36 B <- tbase(x, deg=3, n.knots=4)
R 37 # Estimate coefficients
R 38 beta <- qr.coef(qr(B), y)
R 39
R 40 # Plot fitted function
R 41 plot(x, y, xlab="Age", ylab="log(Count)")
R 42 xx<-seq(min(x), max(x), len=1e3)
R 43 BB <- tbase(xx, deg=3, n.knots=4)
R 44 lines(xx, BB%*%beta)
R 45
R 46 # Estimated derivative
R 47 BBd <-tbase(xx, deg=2, n.knots=4)
R 48 gamma <- beta[-1]*c(1:2, rep(3, 4-1))
R 49 plot(xx, BBd%*%gamma, type="l", xlab="Age", ylab="Derivative")
```

(e) For the local smoother we have that

$$\begin{aligned} \frac{\partial}{\partial x} m(x) &= \frac{\partial}{\partial x} \frac{\sum_{k=1}^n w(x_k - x; h) y_k}{\sum_{k=1}^n w(x_k - x; h)} \\ &= \frac{(\sum_{k=1}^n w'(x_k - x; h) y_k) (\sum_{k=1}^n w(x_k - x; h)) - (\sum_{k=1}^n w(x_k - x; h) y_k) (\sum_{k=1}^n w'(x_k - x; h))}{(\sum_{k=1}^n w(x_k - x; h))^2} \end{aligned}$$

For a Gaussian kernel we have that

$$w(t; h) = \frac{1}{\sqrt{2} * \pi h} \exp(-t^2/2h^2) \quad w'(t; h) = -\frac{t}{\sqrt{2} * \pi h^3} \exp(-t^2/2h^2) = -\frac{t}{h^2} w(t; h),$$

```
R 50 # Set bandwidth
R 51 h <- 5
R 52
R 53 # Compute the (predictive) smoothing matrix S
R 54 S <- dnorm(outer(xx, x, "-"),sd=h)
R 55 S <- S / rowSums(S)
R 56
R 57 # Plot data
R 58 plot(x, y, xlab="Age", ylab="log(Count)")
R 59 lines(xx, S%*%y)
R 60
R 61 # Calculate derivative
R 62 D <- outer(xx, x, "-")
R 63 T <- dnorm(D,sd=h)
R 64 plot(xx, -1/h^2 * ((T*D)%*%y * rowSums(T) - T%*%y * rowSums(T*D))
R 65 / rowSums(T)^2, type="l")
```

7. Model answers not being made available.

8. Model answers not being made available.

9. (a) We continue to use the double-indexing used in \hat{y} . Then $\mathbf{S} = (s_{ij;kl})$ with

$$\begin{aligned} s_{ij;kl} &= \hat{y}_{ij} = \hat{m}(x_i, z_j) = \frac{w(x_k - x_i; h)w(z_l - z_j; h)}{\sum_{\kappa=1}^m \sum_{\lambda=1}^n w(x_{\kappa} - x_i; h)w(z_{\lambda} - z_j; h)} \\ &= \frac{w(x_k - x_i; h)}{\sum_{\kappa=1}^m w(x_{\kappa} - x_i; h)} \frac{w(z_l - z_j; h)}{\sum_{\lambda=1}^n w(z_{\lambda} - z_j; h)}. \end{aligned}$$

(b) This follows directly from part (a) and the definition of the Kronecker product. The first fraction in part (a) comes from the matrix \mathbf{S}_1 and the second fraction comes from \mathbf{S}_2 .

(c) This follows directly from the last property mentioned on page 9 of the assignment sheet. The key advantage of this representation is that much smaller matrices are used, so computations are much more efficient and much larger data sets can be considered.

(d) We can use the following R code.

```
R 1 # Set bandwidth
R 2 h <- 0.5
R 3
R 4 # Construct S1
R 5 S1 <- dnorm(outer(us.easting, us.easting, "-"),sd=h)
R 6 S1 <- S1 / rowSums(S1)
R 7
```



```

R 8 # Construct S2
R 9 S2 <- dnorm(outer(us.northing, us.northing, "-"),sd=h)
R 10 S2 <- S2 / rowSums(S2)
R 11
R 12 # Compute fitted values
R 13 fitted.rain <- S1%*%us.rain%*%t(S2)
R 14
R 15 # Plot results
R 16 library(maps)
R 17 m <- map("state", mar=rep(0.5,4))
R 18 image(us.easting-360, us.northing, fitted.rain, add=TRUE)
R 19 lines(m)

```

10. (a) We have that

$$\begin{aligned}
\mathbf{B} &= \begin{pmatrix} B_{11}(x_1, z_1) & \dots & B_{l_1+r-1,1}(x_1, z_1) & B_{12}(x_1, z_1) & \dots & B_{l_1+r-1,2}(x_1, z_1) & \dots & B_{l_1+r-1,l+2+r-1}(x_1, z_1) \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ B_{11}(x_m, z_1) & \dots & B_{l_1+r-1,1}(x_m, z_1) & B_{12}(x_m, z_1) & \dots & B_{l_1+r-1,2}(x_m, z_1) & \dots & B_{l_1+r-1,l+2+r-1}(x_m, z_1) \\ B_{11}(x_1, z_2) & \dots & B_{l_1+r-1,1}(x_1, z_2) & B_{12}(x_1, z_2) & \dots & B_{l_1+r-1,2}(x_1, z_2) & \dots & B_{l_1+r-1,l+2+r-1}(x_1, z_2) \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ B_{11}(x_m, z_n) & \dots & B_{l_1+r-1,1}(x_m, z_n) & B_{12}(x_m, z_n) & \dots & B_{l_1+r-1,2}(x_m, z_n) & \dots & B_{l_1+r-1,l+2+r-1}(x_m, z_n) \end{pmatrix} \\
&= \begin{pmatrix} B_1^{(1)}(x_1)B_1^{(2)}(z_1) & \dots & B_{l_1+r-1}^{(1)}(x_1)B_1^{(2)}(z_1) & B_1^{(1)}(x_1)B_2^{(2)}(z_1) & \dots & B_{l_1+r-1}^{(1)}(x_1)B_2^{(2)}(z_1) & \dots & B_{l_1+r-1}^{(1)}(x_1)B_{l_2+r-1}^{(2)}(z_1) \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ B_1^{(1)}(x_m)B_1^{(2)}(z_1) & \dots & B_{l_1+r-1}^{(1)}(x_m)B_1^{(2)}(z_1) & B_1^{(1)}(x_m)B_2^{(2)}(z_1) & \dots & B_{l_1+r-1}^{(1)}(x_m)B_2^{(2)}(z_1) & \dots & B_{l_1+r-1}^{(1)}(x_m)B_{l_2+r-1}^{(2)}(z_1) \\ B_1^{(1)}(x_1)B_1^{(2)}(z_2) & \dots & B_{l_1+r-1}^{(1)}(x_1)B_1^{(2)}(z_2) & B_1^{(1)}(x_1)B_2^{(2)}(z_2) & \dots & B_{l_1+r-1}^{(1)}(x_1)B_2^{(2)}(z_2) & \dots & B_{l_1+r-1}^{(1)}(x_1)B_{l_2+r-1}^{(2)}(z_2) \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ B_1^{(1)}(x_m)B_1^{(2)}(z_n) & \dots & B_{l_1+r-1}^{(1)}(x_m)B_1^{(2)}(z_n) & B_1^{(1)}(x_m)B_2^{(2)}(z_n) & \dots & B_{l_1+r-1}^{(1)}(x_m)B_2^{(2)}(z_n) & \dots & B_{l_1+r-1}^{(1)}(x_m)B_{l_2+r-1}^{(2)}(z_n) \end{pmatrix} \\
&= \begin{pmatrix} B_1^{(2)}(z_1) & \dots & B_{l_2+r-1}^{(2)}(z_1) \\ \vdots & \ddots & \vdots \\ B_1^{(2)}(z_n) & \dots & B_{l_2+r-1}^{(2)}(z_n) \end{pmatrix} \otimes \begin{pmatrix} B_1^{(1)}(x_1) & \dots & B_{l_1+r-1}^{(1)}(x_1) \\ \vdots & \ddots & \vdots \\ B_1^{(1)}(x_m) & \dots & B_{l_1+r-1}^{(1)}(x_m) \end{pmatrix} \\
&= \mathbf{s}^{(2)} \otimes \mathbf{s}^{(1)}
\end{aligned}$$

(b) Using that $\mathbf{B} = \mathbf{B}^{(2)} \otimes \mathbf{B}^{(1)}$ we have that

$$\begin{aligned}
\mathbf{S} &= \mathbf{B}(\mathbf{B}^\top \mathbf{B})^{-1} \mathbf{B}^\top \\
&= (\mathbf{B}^{(2)} \otimes \mathbf{B}^{(1)}) \underbrace{\left((\mathbf{B}^{(2)} \otimes \mathbf{B}^{(1)})^\top (\mathbf{B}^{(2)} \otimes \mathbf{B}^{(1)}) \right)^{-1}}_{= (\mathbf{B}^{(2)\top} \otimes \mathbf{B}^{(1)\top}) (\mathbf{B}^{(2)} \otimes \mathbf{B}^{(1)})^{-1} = (\mathbf{B}^{(2)\top} \mathbf{B}^{(2)})^{-1} \otimes (\mathbf{B}^{(1)\top} \mathbf{B}^{(1)})^{-1}} (\mathbf{B}^{(2)} \otimes \mathbf{B}^{(1)})^\top \\
&= \underbrace{\left((\mathbf{B}^{(2)\top} \otimes \mathbf{B}^{(1)\top}) (\mathbf{B}^{(2)} \otimes \mathbf{B}^{(1)}) \right)^{-1}}_{= (\mathbf{B}^{(2)\top} \mathbf{B}^{(2)})^{-1} \otimes (\mathbf{B}^{(1)\top} \mathbf{B}^{(1)})^{-1}} (\mathbf{B}^{(2)} \otimes \mathbf{B}^{(1)})^\top \\
&= (\mathbf{B}^{(2)} (\mathbf{B}^{(2)\top} \mathbf{B}^{(2)})^{-1} \mathbf{B}^{(2)}) \otimes (\mathbf{B}^{(1)} (\mathbf{B}^{(1)\top} \mathbf{B}^{(1)})^{-1} \mathbf{B}^{(1)}) = \mathbf{S}^{(2)} \otimes \mathbf{S}^{(1)}
\end{aligned}$$

The result follows from the last property mentioned on page 9 of the assignment sheet.

(c) For P-splines we have, using $\mathbf{D} = \begin{pmatrix} 1 & -1 & 0 & \dots & 0 & 0 \\ 0 & 1 & -1 & \ddots & 0 & 0 \\ \vdots & \ddots & \ddots & \ddots & \vdots & \vdots \\ 0 & \dots & \dots & \dots & 1 & -1 \end{pmatrix}$

$$\begin{aligned}
\mathbf{S} &= \mathbf{B}(\mathbf{B}^\top \mathbf{B} + \lambda \mathbf{D}^\top \mathbf{D})^{-1} \mathbf{B}^\top \\
&= \mathbf{B}(\mathbf{B}^\top \mathbf{B} + \lambda \mathbf{D}^\top \mathbf{D})^{-1} \mathbf{B}^\top \\
\mathbf{S}^{(2)} \otimes \mathbf{S}^{(1)} &= (\mathbf{B}^{(2)}(\mathbf{B}^{(2)\top} \mathbf{B}^{(2)} + \lambda \mathbf{D}^{(2)\top} \mathbf{D}^{(2)})^{-1} \mathbf{B}^{(2)\top}) \otimes (\mathbf{B}^{(1)}(\mathbf{B}^{(1)\top} \mathbf{B}^{(1)} + \lambda \mathbf{D}^{(1)\top} \mathbf{D}^{(1)})^{-1} \mathbf{B}^{(1)\top}) \\
&= (\mathbf{B}^{(2)} \otimes \mathbf{B}^{(1)}) \left((\mathbf{B}^{(2)\top} \mathbf{B}^{(2)} + \lambda \mathbf{D}^{(2)\top} \mathbf{D}^{(2)}) \otimes (\mathbf{B}^{(1)\top} \mathbf{B}^{(1)} + \lambda \mathbf{D}^{(1)\top} \mathbf{D}^{(1)}) \right)^{-1} (\mathbf{B}^{(2)} \otimes \mathbf{B}^{(1)})^\top \\
&= (\mathbf{B}^{(2)} \otimes \mathbf{B}^{(1)}) \left((\mathbf{B}^{(2)} \otimes \mathbf{B}^{(1)})^\top (\mathbf{B}^{(2)} \otimes \mathbf{B}^{(1)}) + \lambda (\mathbf{D}^{(2)} \otimes \mathbf{B}^{(1)})^\top (\mathbf{D}^{(2)} \otimes \mathbf{B}^{(1)}) \right. \\
&\quad \left. + \lambda (\mathbf{B}^{(2)} \otimes \mathbf{D}^{(1)})^\top (\mathbf{B}^{(2)} \otimes \mathbf{D}^{(1)}) + \lambda^2 (\mathbf{D}^{(2)} \otimes \mathbf{D}^{(1)})^\top (\mathbf{D}^{(2)} \otimes \mathbf{D}^{(1)}) \right)^{-1} (\mathbf{B}^{(2)} \otimes \mathbf{B}^{(1)})^\top \\
&= \mathbf{B}^\top (\mathbf{B}^\top \mathbf{B} + \lambda (\mathbf{D}^{(2)} \otimes \mathbf{B}^{(1)})^\top (\mathbf{D}^{(2)} \otimes \mathbf{B}^{(1)}) + \lambda (\mathbf{B}^{(2)} \otimes \mathbf{D}^{(1)})^\top (\mathbf{B}^{(2)} \otimes \mathbf{D}^{(1)}) \\
&\quad + \lambda^2 (\mathbf{D}^{(2)} \otimes \mathbf{D}^{(1)})^\top (\mathbf{D}^{(2)} \otimes \mathbf{D}^{(1)}))^{-1} \mathbf{B}^\top
\end{aligned}$$

Thus the two approaches are not the same, as they correspond to different penalties.

(d) We can use the following R code.

```

R1 # Construct S1
R2 B1 <- bbase(us.easting, n.knots=30)
R3 S1 <- B1%%solve(crossprod(B1))%%t(B1)
R4
R5 # Construct S2
R6 B2 <- bbase(us.northing, n.knots=7)
R7 S2 <- B2%%solve(crossprod(B2))%%t(B2)
R8
R9
R10 # Compute fitted values
R11 fitted.rain <- S1%%us.rain%%t(S2)
R12
R13 # Plot results
R14 library(maps)
R15 m <- map("state", mar=rep(0.5,4))
R16 image(us.easting-360, us.northing, fitted.rain, add=TRUE)
R17 lines(m)

```

11. (a) Consider a hypothetical second set of observations \mathbf{y}' observed at the same covariates as \mathbf{y} . Then $\hat{\mathbf{y}} = \mathbb{E}(\mathbf{y}'|\mathbf{y})$.

The joint distribution of \mathbf{y} and \mathbf{y}' is

$$\begin{pmatrix} \mathbf{y} \\ \mathbf{y}' \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} \mathbf{0} \\ \mathbf{0} \end{pmatrix}, \begin{pmatrix} \mathbf{K} + \sigma^2 \mathbf{I} & \mathbf{K} \\ \mathbf{K}^\top & \mathbf{K} + \sigma^2 \mathbf{I} \end{pmatrix} \right),$$

and thus

$$\mathbf{y}'|\mathbf{y} \sim \mathcal{N}(\mathbf{K}(\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}, \mathbf{K} - \mathbf{K}(\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{K} + \sigma^2 \mathbf{I})$$

yielding

$$\hat{\mathbf{y}} = \mathbb{E}(\mathbf{y}'|\mathbf{y}) = \mathbf{K}(\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}$$

(b) Using $\mathbf{K} = \mathbf{\Gamma}\mathbf{\Lambda}\mathbf{\Gamma}^\top$

$$\begin{aligned}
\hat{\mathbf{y}} &= \mathbf{K}(\mathbf{K} + \sigma^2\mathbf{I})^{-1}\mathbf{y} = \mathbf{\Gamma}\mathbf{\Lambda}\mathbf{\Gamma}^\top \left(\mathbf{\Gamma}\mathbf{\Lambda}\mathbf{\Gamma}^\top + \sigma^2 \underbrace{\mathbf{I}}_{=\mathbf{\Gamma}\mathbf{\Gamma}^\top} \right)^{-1} \mathbf{y} \\
&= \mathbf{\Gamma}\mathbf{\Lambda}\mathbf{\Gamma}^\top (\mathbf{\Gamma}(\mathbf{\Lambda} + \sigma^2\mathbf{I})\mathbf{\Gamma}^\top)^{-1} \mathbf{y} \\
&= \mathbf{\Gamma}\mathbf{\Lambda} \underbrace{\mathbf{\Gamma}^\top\mathbf{\Gamma}}_{=\mathbf{I}} (\mathbf{\Lambda} + \sigma^2\mathbf{I})^{-1} \mathbf{\Gamma}^\top \mathbf{y} \\
&= \mathbf{\Gamma} \underbrace{\mathbf{\Lambda}(\mathbf{\Lambda} + \sigma^2\mathbf{I})^{-1}}_{=\text{diag}\left(\frac{\lambda_i}{\lambda_i + \sigma^2}\right)} \mathbf{\Gamma} \mathbf{y} \\
&= \mathbf{\Gamma} \text{diag}\left(\frac{\lambda_i}{\lambda_i + \sigma^2}\right) \mathbf{\Gamma}^\top \mathbf{y}
\end{aligned}$$

(c) (i) Let's assume that the values of the covariates (easting and northing) are ordered as follows

$$\begin{aligned}
\mathbf{x}_1 &= (x_1, z_1) \\
\mathbf{x}_2 &= (x_2, z_1) \\
&\vdots \\
\mathbf{x}_m &= (x_m, z_1) \\
\mathbf{x}_{m+1} &= (x_1, z_2) \\
&\vdots \\
\mathbf{x}_{mn} &= (x_m, z_n)
\end{aligned}$$

Then

$$\begin{aligned}
\mathbf{K} &= \begin{bmatrix} k^{(1)}(x_1, x_1)k^{(2)}(z_1, z_1) & k^{(1)}(x_1, x_2)k^{(2)}(z_1, z_1) & \dots & k^{(1)}(x_1, x_m)k^{(2)}(z_1, z_1) & k^{(1)}(x_1, x_1)k^{(2)}(z_1, z_2) & \dots & k^{(1)}(x_1, x_m)k^{(2)}(z_1, z_n) \\ k^{(1)}(x_2, x_1)k^{(2)}(z_1, z_1) & k^{(1)}(x_2, x_2)k^{(2)}(z_1, z_1) & \dots & k^{(1)}(x_2, x_m)k^{(2)}(z_1, z_1) & k^{(1)}(x_2, x_1)k^{(2)}(z_1, z_2) & \dots & k^{(1)}(x_2, x_m)k^{(2)}(z_1, z_n) \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ k^{(1)}(x_m, x_1)k^{(2)}(z_1, z_1) & k^{(1)}(x_m, x_2)k^{(2)}(z_1, z_1) & \dots & k^{(1)}(x_m, x_m)k^{(2)}(z_1, z_1) & k^{(1)}(x_m, x_1)k^{(2)}(z_1, z_2) & \dots & k^{(1)}(x_m, x_m)k^{(2)}(z_1, z_n) \\ k^{(1)}(x_1, x_1)k^{(2)}(z_2, z_1) & k^{(1)}(x_1, x_2)k^{(2)}(z_2, z_1) & \dots & k^{(1)}(x_1, x_m)k^{(2)}(z_2, z_1) & k^{(1)}(x_1, x_1)k^{(2)}(z_2, z_2) & \dots & k^{(1)}(x_1, x_m)k^{(2)}(z_2, z_n) \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ k^{(1)}(x_m, x_1)k^{(2)}(z_n, z_1) & k^{(1)}(x_m, x_2)k^{(2)}(z_n, z_1) & \dots & k^{(1)}(x_m, x_m)k^{(2)}(z_n, z_1) & k^{(1)}(x_m, x_1)k^{(2)}(z_n, z_2) & \dots & k^{(1)}(x_m, x_m)k^{(2)}(z_n, z_n) \end{bmatrix} \\
&= \begin{bmatrix} k^{(2)}(z_1, z_1) & \dots & k^{(2)}(z_1, z_n) \\ \vdots & \ddots & \vdots \\ k^{(2)}(z_n, z_1) & \dots & k^{(2)}(z_n, z_n) \end{bmatrix} \otimes \begin{bmatrix} k^{(1)}(x_1, x_1) & \dots & k^{(1)}(x_1, x_m) \\ \vdots & \ddots & \vdots \\ k^{(1)}(x_m, x_1) & \dots & k^{(1)}(x_m, x_m) \end{bmatrix} \\
&= \mathbf{K}^{(2)} \otimes \mathbf{K}^{(1)}
\end{aligned}$$

(ii) Recall that

$$\hat{\mathbf{y}} = \mathbf{\Gamma}\mathbf{\Lambda}(\mathbf{\Lambda} + \sigma^2\mathbf{I})^{-1}\mathbf{\Gamma}^\top \mathbf{y} = (\mathbf{\Gamma}^{(2)} \otimes \mathbf{\Gamma}^{(1)})(\mathbf{\Lambda}^{(2)} \otimes \mathbf{\Lambda}^{(1)})(\mathbf{\Lambda}^{(2)} \otimes \mathbf{\Lambda}^{(1)} + \sigma^2\mathbf{I})^{-1}(\mathbf{\Gamma}^{(2)} \otimes \mathbf{\Gamma}^{(1)})^\top \mathbf{y}$$

Starting from the right of the expression we will now move from a vector view, which involves Kronecker products, to a ‘‘marginal’’ matrix view. We start with

$$\mathbf{u} = (\mathbf{\Gamma}^{(2)} \otimes \mathbf{\Gamma}^{(1)})^\top \mathbf{y} = \mathbf{\Gamma}^{(2)\top} \otimes \mathbf{\Gamma}^{(1)\top} \mathbf{y},$$

which is a vector which has the same entries as the matrix

$$\mathbf{U} = \mathbf{\Gamma}^{(1)\top} \mathbf{Y} \mathbf{\Gamma}^{(2)}.$$

Furthermore $\mathbf{\Lambda}^{(2)} \otimes \mathbf{\Lambda}^{(1)}$ is a diagonal matrix containing the product of each pair of eigenvalues of $\mathbf{K}^{(2)}$ and $\mathbf{K}^{(1)}$. If we denote these eigenvalues by $\lambda_i^{(1)}$ and $\lambda_j^{(2)}$ then

$$\mathbf{\Lambda}^{(2)} \otimes \mathbf{\Lambda}^{(1)} = \text{diag}(\lambda_j^{(2)}) \otimes \text{diag}(\lambda_i^{(1)}) = \text{diag}(\lambda_j^{(2)} \lambda_i^{(1)})$$

with i being the fast-moving index. Thus

$$(\mathbf{\Lambda}^{(2)} \otimes \mathbf{\Lambda}^{(1)})(\mathbf{\Lambda}^{(2)} \otimes \mathbf{\Lambda}^{(1)} + \sigma^2 \mathbf{I})^{-1} = \text{diag} \left(\frac{\lambda_j^{(2)} \lambda_i^{(1)}}{\lambda_j^{(2)} \lambda_i^{(1)} + \sigma^2} \right),$$

and the vector

$$\mathbf{v} = (\mathbf{\Lambda}^{(2)} \otimes \mathbf{\Lambda}^{(1)})(\mathbf{\Lambda}^{(2)} \otimes \mathbf{\Lambda}^{(1)} + \sigma^2 \mathbf{I})^{-1} \mathbf{u} = \text{diag} \left(\frac{\lambda_j^{(2)} \lambda_i^{(1)}}{\lambda_j^{(2)} \lambda_i^{(1)} + \sigma^2} \right) \mathbf{u}$$

has the same entries as the matrix

$$\mathbf{V} = \mathbf{D} \odot \mathbf{U}$$

where $D_{ij} = \frac{\lambda_j^{(2)} \lambda_i^{(1)}}{\lambda_j^{(2)} \lambda_i^{(1)} + \sigma^2}$ and \odot stands for element-wise multiplication.

Finally, the vector

$$\hat{\mathbf{y}} = (\mathbf{\Gamma}^{(2)} \otimes \mathbf{\Gamma}^{(1)}) \mathbf{v}$$

has the same entries as the matrix

$$\begin{aligned} \hat{\mathbf{Y}} &= \mathbf{\Gamma}^{(1)} \mathbf{V} \mathbf{\Gamma}^{(2)\top} \\ &= \mathbf{\Gamma}^{(1)} (\mathbf{D} \odot \mathbf{U}) \mathbf{\Gamma}^{(2)\top} \\ &= \mathbf{\Gamma}^{(1)} \left(\mathbf{D} \odot \left(\mathbf{\Gamma}^{(1)\top} \mathbf{Y} \mathbf{\Gamma}^{(2)} \right) \right) \mathbf{\Gamma}^{(2)\top} \end{aligned}$$

(d) (i) A naive implementation would be ...

```
R 18 # For simplicity we use an isotropic squared exponential kernel
R 19 x <- expand.grid(easting=us.easting, northing=us.northing)
R 20 # Compute matrix of pairwise distances
R 21 dists <- sqrt(outer(x$easting,x$easting,"-")^2+outer(x$northing,x$northing,"-")^2)
R 22 # Half-way sensible guesses of the hyperparameters (see below for a way of estimating
R 23 sigma2 <- 500
R 24 tau2 <- 6000
R 25 phi <- 0.7
R 26 # Define covariance/kernel matrix K
R 27 library(geoR)
R 28 K <- tau2*matern(dists,phi=phi,kappa=1.5)
R 29 # Naive calculation (slow!)
R 30 y <- c(us.rain)
R 31 yhat <- K%%solve(K+sigma2*diag(nrow(K)),y)
R 32 fitted.rain <- matrix(yhat, nrow=nrow(us.rain), ncol=ncol(us.rain))
R 33 # Plot results
R 34 library(maps)
R 35 m <- map("state", mar=rep(0.5,4))
R 36 image(us.easting-360, us.northing, fitted.rain, add=TRUE)
R 37 lines(m)
```

(ii) Exploiting the properties of the Kronecker product makes the calculations a lot more efficient. We make a slightly different assumption about the kernel function in this part, as we need the kernel function to be separable.

```

R 38 K1 <- matern(abs(outer(us.easting,us.easting,"-")), phi=phi, kappa=1.5)
R 39 K2 <- matern(abs(outer(us.northing,us.northing,"-")), phi=phi, kappa=1.5)
R 40 # We could use K <- K2 %x% K1 and the code from above
R 41 K1.eigen <- eigen(K1)
R 42 K2.eigen <- eigen(K2)
R 43 D0 <- K1.eigen$values%*%t(K2.eigen$values)
R 44 D <- tau2*D0/(tau2*D0+sigma2)
R 45 fitted.rain <- K1.eigen$vectors %*% (D * (t(K1.eigen$vectors)%*%us.rain)%*%K2.eigen$vec
R 46 # Plot results
R 47 library(maps)
R 48 m <- map("state", mar=rep(0.5,4))
R 49 image(us.easting-360, us.northing, fitted.rain, add=TRUE)
R 50 lines(m)

```

This even allows tuning the hyperparameters (using maximum marginal likelihood) in a few seconds.

```

R 51 phis <- exp(seq(-2,2,len=50))
R 52 logliks <- lapply(phis, function(phi) {
R 53   K1 <- matern(abs(outer(us.easting,us.easting,"-")), phi=phi, kappa=15)
R 54   K2 <- matern(abs(outer(us.northing,us.northing,"-")), phi=phi, kappa=15)
R 55   # We could use K <- K2 %x% K1 and the code from above
R 56   K1.eigen <- eigen(K1)
R 57   K2.eigen <- eigen(K2)
R 58   D0 <- K1.eigen$values%*%t(K2.eigen$values)
R 59   loglik <- function(pars) { # Function returns negative loglikelihood
R 60     sigma2 <- pars[1]
R 61     tau2 <- pars[2]
R 62     D <- (tau2*D0+sigma2)
R 63     0.5*sum(log(D))+0.5*sum(us.rain * K1.eigen$vectors %*% (1/D * (t(K1.eigen$vect
R 64   })
R 65   opt <- optim(loglik, par=c(1000,5000))
R 66   loglik <- opt$value
R 67   attr(loglik,"pars") <- opt$par
R 68   loglik
R 69 })
R 70 best.idx <- which.min(logliks)
R 71 phi <- phis[best.idx]
R 72 sigma2 <- attr(logliks[[best.idx]],"pars")[1]
R 73 tau2 <- attr(logliks[[best.idx]],"pars")[2]

```

(iii) Using geor:

```

R 74 # Load geor
R 75 library(geor)
R 76 # Convert the US rain ata to geodata format
R 77 us <- as.geodata(cbind(expand.grid(x=us.easting,y=us.northing), c(us.rain)))
R 78 # Estimate variogram (not needed, mostly for historic reasons)
R 79 vario <- variogram(us, max.dist=25)
R 80 plot(vario)
R 81 # Estimate covariance/kernel parameters with maximum likelihood (slow! - might take se
R 82 # reml.est <- likfit(us, ini.cov.pars=c(var(c(us.rain)),1), lik.method="RML", kappa=1.
R 83 # Add estimates to plot of variogram
R 84 # lines(reml.est, col="blue")
R 85 # Estimate the hyperparameters from the variogram (not ideal, but not as slow as likfi
R 86 vario.est <- variofit(vario, kappa=1.5)
R 87 # Add estimates to plot of variogram

```

```
R 88 | lines(vario.est, col="red")
R 89 | # Perform the kriging
R 90 | kc <- krige.conv(us, locations=us$coords, krige = krige.control(obj.m = vario.est))
R 91 | # Plot the fitted surface
R 92 | image(kc, locations=us$coords)
```